

ANM.



⑬ BUNDESREPUBLIK
DEUTSCHLAND



DEUTSCHES
PATENT- UND
MARKENAMT

⑫ Übersetzung der
europäischen Patentschrift

⑨ EP 1 016 071 B 1

⑩ DE 698 03 202 T 2

⑤ Int. Cl.⁷:
G 10 L 11/02
G 10 L 11/06

DE 698 03 202 T 2

⑳ Deutsches Aktenzeichen:	698 03 202.0
㉑ PCT-Aktenzeichen:	PCT/FR98/01979
㉒ Europäisches Aktenzeichen:	98 943 998.9
㉓ PCT-Veröffentlichungs-Nr.:	WO 99/14737
㉔ PCT-Anmeldetag:	16. 9. 1998
㉕ Veröffentlichungstag der PCT-Anmeldung:	25. 3. 1999
㉖ Erstveröffentlichung durch das EPA:	5. 7. 2000
㉗ Veröffentlichungstag der Patenterteilung beim EPA:	16. 1. 2002
㉘ Veröffentlichungstag im Patentblatt:	29. 8. 2002

③① Unionspriorität:
9711640 18. 09. 1997 FR

⑦③ Patentinhaber:
Matra Nortel Communications, Quimper, FR; Eads
Defence and Security Networks, Montigny Le
Bretonneux, FR

⑦④ Vertreter:
WINTER, BRANDL, FÜRNISS, HÜBNER, RÖSS,
KAISER, POLTE, Partnerschaft, 85354 Freising

⑧④ Benannte Vertragsstaaten:
CH, DE, FI, GB, LI, SE

⑦⑤ Erfinder:
LOCKWOOD, Philip, F-95490 Vaureal, FR; LUBIARZ,
Stephane, F-95520 Osny, FR

⑤④ VERFAHREN UND VORRICHTUNG ZUR SPRACHDETEKTION

Anmerkung: Innerhalb von neun Monaten nach der Bekanntmachung des Hinweises auf die Erteilung des europäischen Patents kann jedermann beim Europäischen Patentamt gegen das erteilte europäische Patent Einspruch einlegen. Der Einspruch ist schriftlich einzureichen und zu begründen. Er gilt erst als eingelegt, wenn die Einspruchsgebühr entrichtet worden ist (Art. 99 (1) Europäisches Patentübereinkommen).

Die Übersetzung ist gemäß Artikel II § 3 Abs. 1 IntPatÜG 1991 vom Patentinhaber eingereicht worden. Sie wurde vom Deutschen Patent- und Markenamt inhaltlich nicht geprüft.

DE 698 03 202 T 2

Best Available Copy

15.04.02

Beschreibung

5 Die vorliegende Erfindung betrifft digitale Verfahren zur Verarbeitung von Sprachsignalen. Sie betrifft insbesondere Verfahren, die eine Erfassung von Stimmaktivität anwenden, um differenzierte Verarbeitungen je nachdem durchzuführen, ob das Signal eine Stimmaktivität aufweist
10 oder nicht.

Die betreffenden digitalen Verfahren beziehen sich auf verschiedene Fachgebiete: Sprachcodierung für die Übertragung oder Speicherung oder Erkennung von Sprache,
15 Verminderung von Rauschen, Echounterdrückung usw.

Die Verfahren zur Erfassung von Stimmaktivität haben als hauptsächliche Schwierigkeit die Unterscheidung zwischen der Stimmaktivität und dem sie begleitenden Lärm bzw. Rauschen. Die Zuhilfenahme eines klassischen Rauschunterdrückungsverfahrens gestattet es nicht, diese Schwierigkeit zu behandeln, da diese Verfahren wiederum Schätzungen des Rauschens anwenden, die von dem Grad der Stimmaktivität des Signals abhängen. Dieses Problem ist
20 beispielsweise in der Schrift US-A-5659622 beschrieben.

Ein Hauptziel der vorliegenden Erfindung ist es, die Robustheit der Verfahren zur Erfassung von Stimmaktivität gegen Rauschen zu verbessern. Um dieses Ziel zu erreichen, wird ein Verfahren gemäß den Angaben in Anspruch 1
30 vorgeschlagen.

So schlägt die Erfindung ein Verfahren zum Erfassen von Stimmaktivität in einem in aufeinanderfolgenden
35 Blöcken behandelten digitalen Sprachsignal vor, bei dem das Sprachsignal einer Rauschunterdrückung unter Berücksichtigung

15.04.02

sichtigung von Schätzungen des im Signal enthaltenen Rauschens unterzogen wird, die für jeden Block auf eine Weise aktualisiert werden, die von zumindest einem für den betreffenden Block bestimmten Grad der Stimmaktivität abhängt. Erfindungsgemäß wird eine apriorische Rauschunterdrückung des Sprachsignals eines jeden Blocks auf der Grundlage von Schätzungen des Rauschens durchgeführt, die bei der Behandlung von mindestens einem vorhergehenden Block erhalten wurden, und die Variationen der Energie des apriorisch rauschunterdrückten Signals analysiert werden, um den Grad der Stimmaktivität des Blocks zu erfassen.

Der Umstand, daß die Erfassung der Stimmaktivität (gemäß einem Verfahren, bei dem es sich im wesentlichen um jegliches bekannte Verfahren handeln kann) auf der Grundlage eines apriorisch rauschunterdrückten Signals durchgeführt wird, verbessert wesentlich die Leistungsfähigkeit dieser Erfassung, wenn das Umgebungsrauschen relativ stark ist.

In der Folge der vorliegenden Beschreibung wird das erfindungsgemäße Verfahren zum Erfassen von Stimmaktivität an einem System zur Rauschunterdrückung eines Sprachsignals veranschaulicht. Es ist ersichtlich, daß dieses Verfahren Anwendungen in zahlreichen weiteren Arten der digitalen Sprachverarbeitung finden kann, bei denen es erwünscht ist, über eine Information bezüglich des Grades der Stimmaktivität des verarbeiteten Signals zur verfügen: Codierung, Erkennung, Echounterdrückung usw.

Weitere Details und Vorteile der vorliegenden Erfindung ergeben sich aus der nachfolgenden Beschreibung von nicht-einschränkenden Ausführungsbeispielen unter Bezugnahme auf die beigefügte Zeichnung.

Es zeigt:

- Fig. 1 eine schematische Übersicht eines Rauschunterdrückungssystems, das die vorliegende Erfindung anwendet;
- Fig. 2 und 3 Organigramme von Prozeduren, die durch einen Stimmaktivitätsdetektor des Systems von Fig. 1 angewendet werden;
- Fig. 4 ein Diagramm, das die Zustände eines Automaten zur Erfassung von Stimmaktivität darstellt;
- Fig. 5 ein Diagramm zur Veranschaulichung der Variationen eines Grades der Stimmaktivität;
- Fig. 6 eine schematische Übersicht eines Moduls zur Überbewertung von Rauschen des Systems von Fig. 1;
- Fig. 7 ein Diagramm zur Veranschaulichung der Berechnung einer Maskierungskurve; und
- Fig. 8 ein Diagramm zur Veranschaulichung der Nutzung der Maskierungskurven in dem System von Fig. 1.

Das in Fig. 1 dargestellte System zur Rauschunterdrückung behandelt ein digitales Sprachsignal s . Ein Fensterbildungsmodul 10 bringt dieses Signal s in die Form von aufeinanderfolgenden Fenstern oder Blöcken, die jeweils aus einer Anzahl N von Abtastproben eines digitalen Signals bestehen. Auf klassische Weise können diese Blöcke gegenseitige Überlappungen aufweisen. In der nachfolgenden Beschreibung wird angenommen, ohne daß dies einschränkend gedacht ist, daß die Blöcke aus $N = 256$ Abtastproben mit einer Abtastrate F_e von 8 kHz bestehen, mit einer Hamming-Wichtung in jedem Fenster, und Überlappungen von 50% zwischen aufeinanderfolgenden Fenstern.

Der Signalblock wird durch ein Modul 11, das einen klassischen Algorithmus der schnellen Fourier-Transformation (TFR) für die Berechnung des Moduls des Spektrums

15.04.02

des Signals anwendet, in den Frequenzbereich transformiert. Das Modul 11 liefert somit eine mit $S_{n,f}$ bezeichnete Gesamtheit von $N = 256$ Frequenzkomponenten des Sprachsignals, wobei n die Nummer des momentanen Blocks und f eine Frequenz des diskreten Spektrums bezeichnet. Aufgrund der Eigenschaften der digitalen Signale im Frequenzbereich werden nur die $N/2 = 128$ ersten Abtastproben verwendet.

Für die Berechnung der Schätzungen des in dem Signal enthaltenen Rauschens wird nicht die am Ausgang der schnellen Fourier-Transformation verfügbare Frequenzauflösung verwendet, sondern eine schwächere Auflösung, die durch eine Anzahl I von Frequenzbändern bestimmt ist, welche das Band $[0, F_e/2]$ des Signals abdeckt. Jedes Band i ($1 \leq i \leq I$) erstreckt sich zwischen einer unteren Frequenz $f(i-1)$ und einer oberen Frequenz $f(i)$, wobei $f(0) = 0$, und $f(I) = F_e/2$. Dieses Zerschneiden in Frequenzbänder kann gleichförmig ($f(i) - f(i-1) = F_e/2I$) sein. Es kann auch nicht gleichförmig sein (z.B. gemäß einer Barks-Skala). Ein Modul 12 berechnet die jeweiligen Mittelwerte der Spektralkomponenten $S_{n,f}$ des Sprachsignals pro Bändern, beispielsweise durch eine gleichförmige Wichtung wie etwa:

$$S_{n,i} = \frac{1}{f(i) - f(i-1)} \sum_{f \in [f(i-1), f(i)]} S_{n,f} \quad (1)$$

Diese Mittelwertbildung vermindert die Schwankungen zwischen den Bändern durch Mitteln der Beiträge des Rauschens in diesen Bändern, wodurch die Varianz des Schätzers des Rauschens vermindert wird. Des weiteren gestattet diese Mittelwertbildung eine starke Verringerung der Komplexität des Systems.

Die gemittelten Spektralkomponenten $S_{n,i}$ werden an ein Modul 15 für die Erfassung von Stimmaktivität und an ein Modul 16 zur Schätzung des Rauschens adressiert. Diese beiden Module 15, 16 arbeiten insofern gemeinsam, als von dem Modul 15 für die verschiedenen Bänder gemessene Stimmaktivitätsgrade $\gamma_{n,i}$ von dem Modul 16 für die Schätzung der Langzeitenergie des Rauschens in den verschiedenen Bändern verwendet werden, während diese Langzeitschätzungen $\hat{B}_{n,i}$ von dem Modul 15 verwendet werden, um eine apriorische Rauschunterdrückung des Sprachsignals in den verschiedenen Bändern vorzunehmen, um die Stimmaktivitätsgrade $\gamma_{n,i}$ zu bestimmen.

Der Betrieb der Module 15 und 16 kann den in Fig. 2 und 3 dargestellten Organigrammen entsprechen.

In den Schritten 17 bis 20 führt das Modul 15 die apriorische Rauschunterdrückung des Sprachsignals in den unterschiedlichen Bändern i für den Signalblock n durch. Diese apriorische Rauschunterdrückung wird gemäß einem klassischen Vorgang zur nichtlinearen Spektralsubtraktion ausgehend von Schätzungen des Rauschens durchgeführt, welche bei einem oder mehreren vorausgegangenen Blöcken erhalten wurden. In Schritt 17 berechnet das Modul 15 mit der Auflösung der Bänder i den Frequenzgang $H_{p,n,i}$ des Filters für die apriorische Rauschunterdrückung gemäß der Formel:

$$H_{p,n,i} = \frac{S_{n,i} - \alpha'_{n-\tau_1,i} \cdot \hat{B}_{n-\tau_1,i}}{S_{n-\tau_2,i}} \quad (2)$$

wobei τ_1 und τ_2 als Anzahl von Blöcken ausgedrückte Verzögerungen sind ($\tau_1 \geq 1$, $\tau_2 \geq 0$), und $\alpha'_{n,i}$ ein Koeffizient der Überbewertung des Rauschens ist, dessen Bestimmung weiter unten erläutert wird. Die Verzögerung τ_1 kann

festgelegt (z.B. $\tau_1 = 1$) oder auch variabel sein. Sie ist umso geringer, je stärker man sich auf die Erfassung der Stimmaktivität verläßt.

- 5 In den Schritten 18 bis 20 werden die Spektralkomponenten $\hat{E}_{p,n,i}$ berechnet gemäß:

$$\hat{E}_{p,n,i} = \max \{ H_{p,n,i} \cdot S_{n,i}, \beta_{p,i} \cdot \hat{E}_{n-\tau_1,i} \} \quad (3)$$

- 10 wobei $\beta_{p,i}$ ein Untergrenzenkoeffizient nahe 0 ist, der klassischerweise dazu dient zu vermeiden, daß das Spektrum des entrauschten Signals negative oder übermäßig schwache Werte annimmt, die ein musikalisches Geräusch hervorrufen würden.

- 15 Die Schritte 17 bis 20 bestehen somit im wesentlichen darin, von dem Spektrum des Signals eine durch den Koeffizienten $\alpha'_{n-\tau_1,i}$ majorierte Schätzung des apriorisch geschätzten Spektrums des Rauschens zu subtrahieren.

- 20 In Schritt 21 berechnet das Modul 15 die Energie des apriorisch rauschunterdrückten Signals in den verschiedenen Bändern i für den Block n : $E_{n,i} = \hat{E}_{p,n,i}^2$. Es berechnet auch einen globalen Mittelwert $E_{n,0}$ der Energie des apriorisch rauschunterdrückten Signals durch eine Summe der Energien pro Band $E_{n,i}$, die mit den Breiten dieser Bänder gewichtet sind. In den nachfolgenden Angaben wird der Index $i = 0$ dazu verwendet, das globale Band des Signals zu bezeichnen.

- 30 In den Schritten 22 und 23 berechnet das Modul 15 für jedes Band i ($0 \leq i \leq I$) eine Größe $\Delta E_{n,i}$, welche für die Kurzzeitvariation der Energie des entrauschten Signals im Band i steht, sowie einen Langzeitwert $\bar{E}_{n,i}$ der Energie des entrauschten Signals im Band i . Die Größe $\Delta E_{n,i}$ kann

berechnet werden durch eine vereinfachte Ableitungsformel: $\Delta E_{n,i} = \left| \frac{E_{n-4,i} + E_{n-3,i} - E_{n-1,i} - E_{n,i}}{10} \right|$. Was die Langzeitenergie $\bar{E}_{n,i}$ betrifft, so kann diese mit Hilfe eines Vergessensfaktors $B1$ wie etwa $0 < B1 < 1$ berechnet werden, nämlich $\bar{E}_{n,i} = B1 \cdot \bar{E}_{n-1,i} + (1 - B1) \cdot E_{n,i}$.

Nach der Berechnung der Energien $E_{n,i}$ des rauschunterdrückten Signals, seiner Kurzzeitvariationen $\Delta E_{n,i}$ und seiner Langzeitwerte $\bar{E}_{n,i}$ auf die in Fig. 2 angegebene Weise berechnet das Modul 15 für jedes Band i ($0 \leq i \leq I$) einen Wert p_i , der für die Evolution der Energie des rauschunterdrückten Signals steht. Diese Berechnung wird in den Schritten 25 bis 36 von Fig. 3 vorgenommen, die für jedes Band i zwischen $i=0$ und $i=I$ durchgeführt werden. Diese Berechnung wendet einen Langzeitschätzer ba_i der Umhüllenden des Rauschens, einen internen Schätzer bi_i und einen Zähler b_i für verrauschte Blöcke an.

In Schritt 25 wird die Größe $\Delta E_{n,i}$ mit einem Schwellwert ϵ_1 verglichen. Wenn der Schwellwert ϵ_1 nicht erreicht wird, wird der Zähler b_i in Schritt 26 um eine Einheit inkrementiert. In Schritt 27 wird der Langzeitschätzer ba_i mit dem Wert der geglätteten Energie $\bar{E}_{n,i}$ verglichen. Falls $ba_i \geq \bar{E}_{n,i}$, wird der Schätzer ba_i gleich dem geglätteten Wert $\bar{E}_{n,i}$ in Schritt 26 genommen, und der Zähler b_i wird auf Null zurückgesetzt. Die Größe p_i , die gleich dem Verhältnis $ba_i / \bar{E}_{n,i}$ genommen wird (Schritt 36), ist somit gleich 1.

Wenn Schritt 27 ergibt, daß $ba_i < \bar{E}_{n,i}$, wird der Zähler b_i in Schritt 29 mit einem Grenzwert b_{max} verglichen. Falls $b_i > b_{max}$, wird angenommen, daß das Signal zu stationär ist, um Stimmaktivität zu unterstützen. Daraufhin wird der oben genannte Schritt 28 durchgeführt, der in der Annahme besteht, daß der Block nur Rauschen

beinhaltet. Falls $b_i \leq b_{\max}$ in Schritt 29, wird der interne Schätzer bi_i in Schritt 33 berechnet gemäß:

$$bi_i = (1-B_m) \cdot \bar{E}_{n,i} + B_m \cdot ba_i \quad (4)$$

5

In dieser Formel steht B_m für einen zwischen 0,90 und 1 liegenden Aktualisierungskoeffizienten. Sein Wert ist je nach dem Zustand eines Automaten für die Erfassung von Stimmaktivität verschieden (Schritte 30 bis 32). Dieser

10 Zustand δ_{n-1} ist derjenige, der bei der Verarbeitung des vorherigen Blockes bestimmt wurde. Falls sich der Automat in einem Zustand der Erfassung von Sprache befindet ($\delta_{n-1} = 2$ in Schritt 30), nimmt der Koeffizient B_m einen Wert B_{mp} an, der sehr nahe bei 1 liegt, damit der Schätzer des

15 Rauschens bei Vorhandensein von Sprache sehr geringfügig aktualisiert wird. Im entgegengesetzten Fall nimmt der Koeffizient B_m einen geringeren Wert B_{ms} an, um in einer Stillephase eine bedeutendere Aktualisierung des Schätzers des Rauschens zu ermöglichen. In Schritt 34 wird der

20 Abstand $ba_i - bi_i$ zwischen dem Langzeitschätzer und dem internen Schätzer des Rauschens mit einem Schwellwert ϵ_2 verglichen. Wenn der Schwellwert ϵ_2 nicht erreicht wird, wird der Langzeitschätzer ba_i in Schritt 35 mit dem Wert des internen Schätzers bi_i aktualisiert. Andernfalls

25 bleibt der Langzeitschätzer ba_i unverändert. Es wird somit vermieden, daß abrupte Variationen aufgrund eines Sprachsignals zu einer Aktualisierung des Schätzers des Rauschens führen.

30 Nach dem Erhalt der Größen p_i nimmt das Modul 15 die Entscheidungen der Stimmaktivität in Schritt 37 vor. Das Modul 15 aktualisiert zuerst den Zustand des Erfassungsautomaten gemäß der für die Gesamtheit des Bandes des Signals berechneten Größe p_0 . Der neue Zustand δ_n des

Automaten hängt von dem vorhergegangenen Zustand δ_{n-1} und von p_0 ab, wie in Fig. 4 dargestellt ist.

Vier Zustände sind möglich: $\delta = 0$ erfaßt Stille bzw.
 5 Abwesenheit von Sprache; $\delta = 2$ erfaßt das Vorhandensein einer Stimmaktivität; und die Zustände $\delta = 1$ und $\delta = 3$ sind dazwischenliegende Zustände des Anstiegs und Abfallens. Wenn sich der Automat im Zustand von Stille ($\delta_{n-1} = 0$) befindet, bleibt er dort, wenn p_0 nicht eine erste
 10 Schwelle SE1 übersteigt, und geht im entgegengesetzten Fall in den Anstiegzustand über. Im Anstiegzustand ($\delta_{n-1} = 1$) kehrt er in den Zustand von Stille zurück, wenn p_0 kleiner als der Schwellwert SE1 ist, geht in den Zustand der Sprache über, wenn p_0 größer als eine über der
 15 Schwelle SE1 liegende Schwelle SE2 ist, und bleibt im Anstiegzustand, falls $SE1 \leq p_0 \leq SE2$. Wenn sich der Automat im Zustand der Sprache ($\delta_{n-1} = 2$) befindet, so bleibt er dort, falls p_0 eine unter der Schwelle SE2 liegende dritte Schwelle SE3 ist, und geht im entgegengesetzten Fall in den Abstiegzustand über. Im Abstieg-
 20 zustand ($\delta_{n-1} = 3$) kehrt der Automat in den Zustand der Sprache zurück, falls p_0 größer als der Schwellwert SE2 ist, kehrt in den Zustand der Stille zurück, wenn p_0 diesseits eines unter dem Schwellwert SE2 liegenden vier-
 25 ten Schwellwerts SE4 ist, und bleibt im Abstiegzustand, falls $SE4 \leq p_0 \leq SE2$.

In Schritt 37 berechnet das Modul 15 des weiteren die Stimmaktivitätsgrade $\gamma_{n,i}$ in jedem Band $i \geq 1$. Dieser
 30 Grad $\gamma_{n,i}$ ist vorzugsweise ein nicht-binärer Parameter, d.h. die Funktion $\gamma_{n,i} = g(p_i)$ ist eine Funktion, die in Abhängigkeit von den durch die Größe p_i angenommenen Werten kontinuierlich zwischen 0 und 1 variiert. Diese Funktion besitzt beispielsweise den in Fig. 5 dargestellten
 35 Verlauf..

Das Modul 16 berechnet die Schätzungen des Rauschens pro Band, die im Rauschunterdrückungsvorgang verwendet werden, unter Anwendung der aufeinanderfolgenden Werte der Komponenten $S_{n,i}$ und der Stimmaktivitätsgrade $\gamma_{n,i}$. Dies entspricht den Schritten 40 bis 42 von Fig. 3. In Schritt 40 wird bestimmt, ob der Automat für die Erfassung von Stimmaktivität aus dem Anstiegszustand in den Zustand der Sprache übergegangen ist. Falls ja, werden die vorausgehend für jedes Band $i \geq 1$ berechneten beiden letzten Schätzungen $\hat{B}_{n-1,i}$ und $\hat{B}_{n-2,i}$ gemäß dem vorausgegangenen Schätzwert $\hat{B}_{n-3,i}$ korrigiert. Diese Korrektur wird durchgeführt, um den Umstand zu berücksichtigen, daß in der Anstiegsphase ($\delta = 1$) die Langzeitschätzungen der Energie des Rauschens in dem Vorgang für die Erfassung von Stimmaktivität (Schritte 30 bis 33) so berechnet werden konnten, als ob das Signal nur Rauschen beinhaltete ($B_m = B_{ms}$), so daß die Gefahr besteht, daß sie mit einem Fehler behaftet sind.

In Schritt 42 aktualisiert das Modul 16 die Schätzungen des Rauschens pro Band gemäß den Formeln:

$$\bar{B}_{n,i} = \lambda_B \cdot \hat{B}_{n-1,i} + (1 - \lambda_B) \cdot S_{n,i} \quad (5)$$

$$\hat{B}_{n,i} = \gamma_{n,i} \cdot \hat{B}_{n-1,i} + (1 - \gamma_{n,i}) \cdot \bar{B}_{n,i} \quad (6)$$

wobei λ_B einen Vergessensfaktor wie etwa $0 < \lambda_B < 1$ bezeichnet. Formel (6) zeigt die Berücksichtigung des nicht-binären Stimmaktivitätsgrades $\gamma_{n,i}$.

Wie obenstehend angegeben wurde, sind die Langzeitschätzungen des Rauschens $\hat{B}_{n,i}$ Gegenstand einer Überbewertung durch ein Modul 45 (Fig. 1), bevor die Rauschunterdrückung mittels nichtlinearer Spektralsubtraktion

vorgenommen wird. Das Modul 45 berechnet den oben genannten Koeffizienten der Überbewertung $\alpha'_{n,i}$ sowie eine majorierte Schätzung $\hat{B}'_{n,i}$ die im wesentlichen $\alpha'_{n,i} \cdot \hat{B}_{n,i}$ entspricht.

5

Die Strukturierung des Überbewertungsmoduls 45 ist in Fig. 6 dargestellt. Die majorierte Schätzung $\hat{B}'_{n,i}$ wird erhalten durch Kombinieren der Langzeitschätzung $\hat{B}_{n,i}$ und eines Maßes $\Delta B_{n,i}^{\max}$ der Veränderlichkeit der Rauschkomponente in dem Band i um seine Langzeitschätzung. Bei dem betrachteten Beispiel ist dieses Kombinieren im wesentlichen eine einfache Summe, die von einem Addierer 46 erstellt wird. Es könnte sich hierbei auch um eine gewichtete Summe handeln.

10

15

Der Überbewertungskoeffizient $\alpha'_{n,i}$ ist gleich dem Verhältnis zwischen der vom Addierer 46 gelieferten Summe $\hat{B}_{n,i} + \Delta B_{n,i}^{\max}$ und der verzögerten Langzeitschätzung $\hat{B}_{n-\tau_3,i}$ (Teiler 47), die nach oben hin durch einen Grenzwert α_{\max} beschränkt ist, beispielsweise $\alpha_{\max} = 4$ (Block 48). Die Verzögerung τ_3 dient gegebenenfalls dazu, in den Anstiegsphasen ($\delta = 1$) den Wert des Überbewertungskoeffizienten $\alpha'_{n,i}$ zu korrigieren, bevor die Langzeitschätzungen durch die Schritte 40 und 41 von Fig. 3 korrigiert worden sind (z.B. $\tau_3 = 3$).

20

25

Die majorierte Schätzung $\hat{B}'_{n,i}$ wird schließlich gleich $\alpha'_{n,i} \cdot \hat{B}_{n-\tau_3,i}$ genommen (Multiplizierer 49).

30

35

Das Maß $\Delta B_{n,i}^{\max}$ der Veränderlichkeit des Rauschens reflektiert die Varianz des Schätzers des Rauschens. Es wird in Abhängigkeit von den Werten von $S_{n,i}$ und von $\hat{B}_{n,i}$ für eine bestimmte Anzahl von vorherigen Blöcken berechnet, an denen das Sprachsignal keine Stimmaktivität in dem Band i aufweist. Es ist eine Funktion der für eine

15.04.02

Anzahl K von Blöcken mit Stille ($n-k \leq n$) berechneten Abstände $|S_{n-k,i} - \hat{B}_{n-k,i}|$. In dem dargestellten Beispiel ist diese Funktion einfach das Maximum (Block 50). Für jeden Block n wird der Grad der Stimmaktivität $\gamma_{n,i}$ mit einem Schwellwert (Block 51) verglichen, um zu entscheiden, ob der in 52-53 berechnete Abstand $|S_{n,i} - \hat{B}_{n,i}|$ in eine Warteschlange 54 mit K Stellen geladen werden muß, die im Ersteingang/Erstausgang-Modus (FIFO) organisiert ist. Falls $\gamma_{n,i}$ den Schwellwert nicht übersteigt (der gleich 0 sein kann, falls die Funktion $g()$ die Form von Fig. 5 besitzt), wird die FIFO nicht versorgt, während sie es im entgegengesetzten Fall wird. Der in der FIFO-54 enthaltene Maximalwert wird dann als Maß $\Delta B_{n,i}^{\max}$ der Veränderlichkeit geliefert.

Das Maß $\Delta B_{n,i}^{\max}$ der Veränderlichkeit kann als Variante in Abhängigkeit von den Werten $S_{n,f}$ (anstatt $S_{n,i}$) und $\hat{B}_{n,i}$ erhalten werden. Anschließend wird auf die gleiche Weise, mit der Ausnahme, daß die FIFO 54 $|S_{n-k,i} - \hat{B}_{n-k,i}|$ nicht enthält, vorgegangen, jedoch eher $\max_{f \in \{f(i-1), f(i)\}} |S_{n-k,f} - \hat{B}_{n-k,i}|$.

Aufgrund der unabhängigen Langzeitschätzungen der Schwankungen des Rauschens $\hat{B}_{n,i}$ und seiner Kurzzeitveränderlichkeit $\Delta B_{n,i}^{\max}$ stellt der majorierte Schätzer $\hat{B}_{n,i}$ eine ausgezeichnete Robustheit des Rauschunterdrückungsverfahrens gegen musikalische Geräusche zur Verfügung.

Eine erste Phase der spektralen Subtraktion wird durch das in Fig. 1 dargestellte Modul 55 verwirklicht. Diese Phase liefert vor der Auflösung der Bänder i ($1 \leq i \leq I$) den Frequenzgang $H_{n,i}^1$ eines ersten Rauschunterdrückungsfilters in Abhängigkeit von den Komponenten $S_{n,i}$ und $\hat{B}_{n,i}$ und den Überbewertungskoeffizienten $\alpha'_{n,i}$. Diese

15.04.00

Berechnung kann für jedes Band i durchgeführt werden gemäß der Formel:

$$H_{n,i}^1 = \frac{\max\{S_{n,i} - \alpha'_{n,i} \cdot \hat{B}_{n,i}, \beta_i^1 \cdot \hat{B}_{n,i}\}}{S_{n-\tau_4,i}} \quad (7)$$

5 wobei τ_4 eine als $\tau_4 \geq 0$ (z.B. $\tau_4 = 0$) bestimmte ganzzahlige Verzögerung ist. In dem Ausdruck (7) stellt der Koeffizient β_i^1 wie der Koeffizient β_{p_i} der Formel (3) eine Untergrenze dar, die klassischerweise zur Vermeidung
10 von negativen oder zu kleinen Werten des rauschunterdrückten Signals dient.

Auf bekannte Weise (EP 0 534 837) könnte der Überbewertungskoeffizient $\alpha'_{n,i}$ in der Formel (7) durch einen
15 anderen Koeffizienten ersetzt werden, der gleich einer Funktion von $\alpha'_{n,i}$ und einer Schätzung des Rauschabstandes (z.B. $S_{n,i}/\hat{B}_{n,i}$) ist, wobei diese Funktion gemäß dem Schätzwert des Rauschabstandes abnehmend ist. Diese Funktion ist somit gleich $\alpha'_{n,i}$ für die kleinsten Werte des
20 Rauschabstandes. Wenn das Signal stark verrauscht ist, ist es nämlich a priori nicht sinnvoll, den Überbewertungsfaktor zu vermindern. Vorteilhaft nimmt diese Funktion für die höchsten Werte des Rauschabstandes gegen Null hin ab. Dies ermöglicht einen Schutz der energie-
25 reichsten Zonen des Spektrums, in denen das Sprachsignal am bedeutendsten ist, wobei die von dem Signal zu subtrahierende Größe somit gegen Null tendiert.

Diese Strategie kann verfeinert werden, indem sie
30 selektiv auf die Harmonischen der Tonfrequenz ("pitch") des Sprachsignals angewendet wird, wenn dieses eine Stimmaktivität aufweist.

Somit wird bei der in Fig. 1 dargestellten Ausführungsform eine zweite Phase der Rauschunterdrückung durch ein Modul 56 zum Schutz der Harmonischen durchgeführt. Dieses Modul berechnet mit der Auflösung der Fourier-Transformierung den Frequenzgang $H_{n,f}^2$ eines zweiten Rauschunterdrückungsfilters in Abhängigkeit von den Parametern $H_{n,i}^1$, $\alpha'_{n,i}$, $\hat{B}_{n,i}$, δ_n , $S_{n,i}$ und der außerhalb der Stillephasen durch ein Modul für die harmonische Analyse 57 berechneten Tonfrequenz $f_p = F_e/T_p$. In einer Stillephase ($\delta_n = 0$) ist das Modul 56 nicht in Betrieb, d.h. $H_{n,f}^2 = H_{n,i}^1$ für jede Frequenz f eines Bandes i . Das Modul 57 kann jegliches bekannte Verfahren für die Analyse des Sprachsignals des Blocks anwenden, um die Periode T_p zu bestimmen, die als ganze Zahl oder Bruchteil von Abtastproben angegeben wird, z.B. ein lineares Prädiktionsverfahren.

Der durch das Modul 56 zur Verfügung gestellte Schutz kann darin bestehen, daß für jede zu einem Band i gehörige Frequenz f durchgeführt wird:

$$H_{n,f}^2 = 1 \text{ falls } \begin{cases} S_{n,i} - \alpha'_{n,i} \cdot \hat{B}_{n,i} > \beta_i^2 \cdot \hat{B}_{n,i} \\ \text{und ganzzahliges } \exists \eta / |f - \eta \cdot f_p| \leq \Delta f / 2 \end{cases} \quad (8)$$

$$\text{andernfalls } H_{n,f}^2 = H_{n,i}^1 \quad (9)$$

$\Delta f = F_e/N$ stellt die spektrale Auflösung der Fourier-Transformation dar. Wenn $H_{n,f}^2 = 1$, ist die von der Komponente $S_{n,f}$ zu substrahierende Größe Null. In dieser Berechnung drücken die Untergrenzenkoeffizienten β_i^2 (z.B. $\beta_i^2 = \beta_i^1$) den Umstand aus, daß bestimmte Harmonische der Tonfrequenz f_p von Rauschen maskiert sein können, so daß es nicht sinnvoll ist, sie zu schützen.

Diese Schutzstrategie wird vorzugsweise für jede der Frequenzen angewendet, die am nächsten zu den Harmonischen von f_p sind, d.h. auf jedes ganzzahlige η .

5 Wenn man mit δf_p die Frequenzauflösung bezeichnet, bei der das Analysemodul 57 die geschätzte Tonfrequenz f_p erzeugt, d.h. daß die reelle Tonfrequenz zwischen $f_p - \delta f_p/2$ und $f_p + \delta f_p/2$ liegt, dann kann der Abstand zwischen der η -ten Harmonischen der reellen Tonfrequenz und ihrer
10 Schätzung $\eta \times f_p$ (Bedingung (9)) bis $\eta \times \delta f_p/2$ gehen. Bei hohen Werten von η kann dieser Abstand größer als die halbe spektrale Auflösung $\Delta f/2$ der Fourier-Transformierten sein. Um diese Unsicherheit zu berücksichtigen und einen guten Schutz der Harmonischen der reellen Tonfre-
15 quenz zu gewährleisten, kann jede der Frequenzen des Intervalls $[\eta \times f_p - \eta \times \delta f_p / 2, \eta \times f_p + \eta \times \delta f_p / 2]$ geschützt werden, d.h. die obenstehende Bedingung (9) kann ersetzt werden durch:

20 ganzzahliges $\exists \eta / |f - \eta \cdot f_p| \leq (\eta \cdot \delta f_p + \Delta f) / 2$ (9')

Diese Schutzart (Bedingung 9') ist von besonderem Interesse, wenn die Werte von η groß sein können, insbesondere falls das Verfahren in einem Breitbandsystem ver-
25 wendet wird.

Für jede geschützte Frequenz kann der korrigierte Frequenzgang $H_{n,f}^2$ gemäß der obenstehenden Angabe gleich 1 sein, was der Subtraktion einer Größe Null im Rahmen der
30 spektralen Subtraktion entspricht, d.h. einem kompletten Schutz der betreffenden Frequenz. Allgemeiner gesagt, dieser korrigierte Frequenzgang $H_{n,f}^2$ könnte je nach dem gewünschten Schutzgrad gleich einem zwischen 1 und $H_{n,f}^1$ liegenden Wert genommen werden, was der Subtraktion einer
35 Größe entspricht, die kleiner als diejenige ist, die zu

15.04.02

subtrahieren wäre, wenn die betreffende Frequenz nicht geschützt wäre.

Die Spektralkomponenten $S_{n,f}^2$ eines rauschunterdrückten Signals werden durch einen Multiplizierer 58 berechnet:

$$S_{n,f}^2 = H_{n,f}^2 \cdot S_{n,f} \quad (10)$$

Dieses Signal $S_{n,f}^2$ wird an ein Modul 60 geliefert, das für jeden Block n eine Maskierungskurve berechnet durch Anwenden eines psychoakustischen Modells der Gehörwahrnehmung durch das menschliche Ohr.

Das Phänomen der Maskierung ist ein von der Funktion des menschlichen Ohrs her bekanntes Prinzip. Wenn zwei Frequenzen gleichzeitig gehört werden, ist es möglich, daß eine von den beiden nicht mehr hörbar ist. Man sagt dann, daß diese maskiert ist.

Es gibt verschiedene Verfahrensweisen für die Berechnung der Maskierungskurven. Beispielsweise kann die von J.D. Johnston ("Transform Coding of Audio Signals Using Perceptual Noise Criteria", IEEE Journal on Selected Area in Communications, Vol. 6, Nr. 2, Februar 1988) entwickelte angewendet werden. Bei dieser Verfahrensweise wird in der Frequenzskala der Barks gearbeitet. Die Maskierungskurve wird als die Faltung der Funktion der spektralen Dehnung der Basilarmembran im Bark-Bereich mit dem anregenden Signal betrachtet, bestehend in der vorliegenden Anwendung aus dem Signal $S_{n,f}^2$. Die spektrale Dehnungsfunktion kann auf die in Fig. 7 dargestellte Weise modelliert werden. Für jedes Bark-Band wird der Beitrag der in Betracht gezogenen niederen und hohen Bänder durch die Funktion der Dehnung der Basilarmembran berechnet:

$$C_{n,q} = \sum_{q'=0}^{q-1} \frac{S_{n,q'}^2}{(10^{10/10})^{(q-q')}} + \sum_{q'=q+1}^Q \frac{S_{n,q'}^2}{(10^{25/10})^{(q'-q)}} \quad (11)$$

wobei die Indices q und q' die Bark-Bänder ($0 \leq q, q' \leq Q$) bezeichnen, und $S_{n,q}^2$ für den Mittelwert der Komponenten $S_{n,f}^2$ des rauschunterdrückten Anregungssignals für die diskreten Frequenzen f steht, die zum Bark-Band q gehören.

10 Der Maskierungsschwellwert $M_{n,q}$ wird erhalten durch das Modul 60 für Bark-Band q gemäß der Formel:

$$M_{n,q} = C_{n,q} / R_q \quad (12)$$

15 in der R_q von dem mehr oder minder stimmhaften Charakter des Signals abhängt. Auf bekannte Weise ist eine mögliche Form von R_q :

$$10 \cdot \log_{10}(R_q) = (A+q) \cdot \chi + B \cdot (1 - \chi) \quad (13)$$

20 wobei $A = 14,5$ und $B = 5,5$. χ bezeichnet einen Stimmhaftigkeitsgrad des Sprachsignals, der zwischen Null (keine Stimmhaftigkeit) und 1 (stark stimmhaftes Signal) variiert. Der Parameter χ kann die bekannte Form aufweisen:

$$\chi = \min \left\{ \frac{SFM}{SFM_{\max}}, 1 \right\} \quad (12)$$

30 wobei SFM in Dezibel das Verhältnis zwischen dem arithmetischen Mittel und dem geometrischen Mittel der Energie der Bark-Bänder angibt, und $SFM_{\max} = -60$ dB.

15.04.02

Das Rauschunterdrückungssystem weist darüber hinaus ein Modul 62 auf, das den Frequenzgang des Rauschunterdrückungsfilters in Abhängigkeit von der durch das Modul 60 berechneten Maskierungskurve $M_{n,q}$ und den durch das Modul 45 berechneten majorierten Schätzungen korrigiert. Das Modul 62 entscheidet über das Rauschunterdrückungsniveau, das tatsächlich erzielt werden soll.

Durch einen Vergleich der Umhüllenden der majorierten Schätzung des Rauschens mit der durch die Maskierungsschwellwerte $M_{n,q}$ gebildeten Umhüllenden wird entschieden, das Signal nur in dem Maße zu entrauschen, in dem die majorierte Schätzung $\hat{B}'_{n,i}$ die Maskierungskurve übersteigt. Dies vermeidet eine nutzlose Unterdrückung von durch Sprache maskiertem Rauschen.

Die neue Antwort $H^3_{n,f}$ für eine zu dem Band i gehörende Frequenz f , die durch das Modul 12 und im Bark-Band q definiert wird, hängt somit von dem relativen Abstand zwischen der majorierten Schätzung $\hat{B}'_{n,i}$ der entsprechenden Spektralkomponente des Rauschens und der Maskierungskurve $M_{n,q}$ folgendermaßen ab:

$$H^3_{n,f} = 1 - (1 - H^2_{n,f}) \cdot \max \left\{ \frac{\hat{B}'_{n,i} - M_{n,q}}{\hat{B}'_{n,i}}, 0 \right\} \quad (14)$$

25

Anders ausgedrückt, die bei dem Vorgang der Spektralsubtraktion mit dem Frequenzgang $H^3_{n,f}$ von einer Spektralkomponente $S_{n,f}$ zu subtrahierende Größe ist im wesentlichen gleich dem Minimum zwischen der bei dem Vorgang der Spektralsubtraktion mit dem Frequenzgang $H^2_{n,f}$ von dieser Spektralkomponente zu subtrahierenden Größe einerseits und dem Bruchteil der majorierten Schätzung $\hat{B}'_{n,i}$ der entsprechenden Spektralkomponente des Rauschens anderer-

30

seits, die gegebenenfalls die Maskierungskurve $M_{n,q}$ übersteigt.

Fig. 8 veranschaulicht das Prinzip der durch das Modul 62 angewendeten Korrektur. Sie zeigt schematisch ein Beispiel für eine auf der Grundlage der Spektralkomponenten $S_{n,f}^2$ des rauschunterdrückten Signals sowie der majorierten Schätzung $\hat{B}'_{n,i}$ des Spektrums des Rauschens berechnete Maskierungskurve $M_{n,q}$. Die schließlich von den Komponenten $S_{n,f}$ zu subtrahierende Größe ist die durch die schraffierten Bereiche dargestellte, d.h. diejenige, die auf den Bruchteil der majorierten Schätzung $\hat{B}'_{n,i}$ der Spektralkomponenten des Rauschens, das die Maskierungskurve übersteigt, begrenzt ist.

Diese Subtraktion wird durch Multiplizieren des Frequenzgangs $H_{n,f}^3$ des Rauschunterdrückungsfilters mit den Spektralkomponenten $S_{n,f}$ des Sprachsignals (Multiplizierer 64) durchgeführt. Ein Modul 65 rekonstruiert sodann das rauschunterdrückte Signal im Zeitbereich mittels Durchführung der schnellen inversen Fourier-Transformierung (TFRI) der vom Multiplizierer 64 gelieferten Abtastproben der Frequenz $S_{n,f}^3$. Bei jedem Block werden einzig die $N/2 = 128$ ersten Abtastproben des durch das Modul 65 erzeugten Signals als endgültiges rauschunterdrücktes Signal s^3 geliefert, nach Rekonstruktion mittels Addition-Überlappung mit den $N/2 = 128$ letzten Abtastproben des vorangegangenen Blocks (Modul 66).

15.04.02

Ansprüche

- 5 1. Verfahren zum Erfassen von Stimmaktivität in einem in
aufeinanderfolgenden Blöcken behandelten digitalen
Sprachsignal (s), bei dem das Sprachsignal einer
Rauschunterdrückung unter Berücksichtigung von Schät-
10 zungen des im Signal enthaltenen Rauschens unterzogen
wird, die für jeden Block auf eine Weise aktualisiert
werden, die von zumindest einem für den betreffenden
Block bestimmten Grad der Stimmaktivität ($\gamma_{n,i}$)
abhängt, dadurch gekennzeichnet, daß eine apriorische
15 Rauschunterdrückung des Sprachsignals eines jeden
Blocks auf der Grundlage von Schätzungen des Rau-
schens durchgeführt wird, die bei der Behandlung von
mindestens einem vorhergehenden Block erhalten wur-
den, und die Variationen der Energie des apriorisch
20 rauschunterdrückten Signals analysiert werden, um den
Grad der Stimmaktivität des Blocks zu erfassen.
2. Verfahren nach Anspruch 1, bei dem der Grad der
Stimmaktivität ($\gamma_{n,i}$) ein nicht-binärer Parameter
ist.
- 25 3. Verfahren nach Anspruch 2, bei dem der Grad der
Stimmaktivität ($\gamma_{n,i}$) eine ständig zwischen 0 und 1
variierende Funktion ist.
- 30 4. Verfahren nach einem der vorhergehenden Ansprüche,
bei dem die Schätzungen des Rauschens in verschiede-
nen Frequenzbändern des Signals erhalten werden, die
apriorische Rauschunterdrückung Band für Band durch-
geführt wird, und ein Grad der Stimmaktivität ($\gamma_{n,i}$)
35 für jedes Band bestimmt wird.

5. Verfahren nach einem der vorhergehenden Ansprüche, bei dem eine Schätzung des Rauschens $\hat{B}_{n,i}$ für den Block n in einem Frequenzband i in der Form

$$\hat{B}_{n,i} = \gamma_{n,i} \cdot \hat{B}_{n-1,i} + (1 - \gamma_{n,i}) \cdot \bar{B}_{n,i}$$

$$\text{mit } \bar{B}_{n,i} = \lambda_B \cdot \hat{B}_{n-1,i} + (1 - \lambda_B) \cdot S_{n,i}$$

erhalten wird,

wobei λ_B ein zwischen 0 und 1 liegender Vergessensfaktor ist, $\gamma_{n,i}$ der für den Block n im Frequenzband i bestimmte Grad der Stimmaktivität ist, und $S_{n,i}$ ein Mittelwert der Amplitude des Spektrums des Sprachsignals des Blocks n im Band i ist.

6. Verfahren nach Anspruch 5, bei dem das apriorisch rauschunterdrückte Signal $\hat{E}_{p,n,i}$ bezüglich eines Blocks n und eines Frequenzbandes i die Form aufweist:

$$\hat{E}_{p,n,i} = \max\{H_{p,n,i} \cdot S_{n,i}, \beta_{p,i} \cdot \hat{B}_{n-\tau_1,i}\}$$

wobei $H_{p,n,i} = \frac{S_{n,i} - \alpha'_{n-\tau_1,i} \cdot \hat{B}_{n-\tau_1,i}}{S_{n-\tau_2,i}}$, τ_1 eine ganze

Zahl von mindestens gleich 1 ist, τ_2 eine ganze Zahl von mindestens gleich 0 ist, $\alpha'_{n-\tau_1,i}$ ein für den Block n- τ_1 und das Band i bestimmter Überbewertungskoeffizient ist, und $\beta_{p,i}$ ein positiver Koeffizient ist.

7. Verfahren nach einem der vorhergehenden Ansprüche, bei dem eine Langzeitschätzung ($\bar{E}_{n,i}$) der Energie des apriorisch rauschunterdrückten Signals ($\hat{E}_{p,n,i}$)

15.04.02

berechnet wird und diese Langzeitschätzung mit einer an dem betreffenden Block berechneten, momentanen Schätzung (b_a) dieser Energie verglichen wird, um den Grad der Stimmaktivität ($\gamma_{n,i}$) des Blocks zu erhalten.

5

8. Stimmaktivität-Erfassungseinrichtung mit einer zum Durchführen eines Verfahrens nach einem der vorhergehenden Ansprüche konzipierten Behandlungseinrichtung.

10

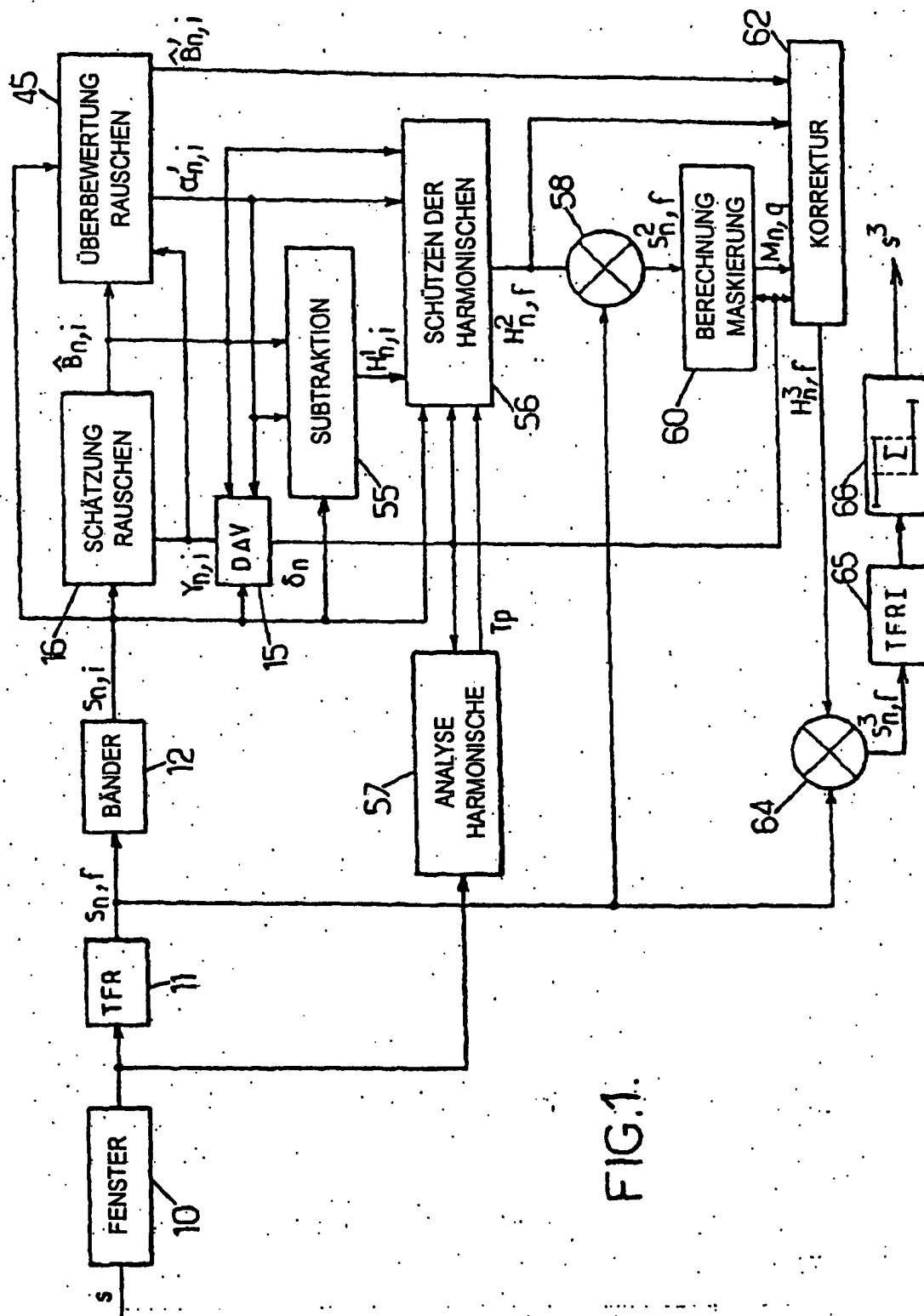


FIG.1.

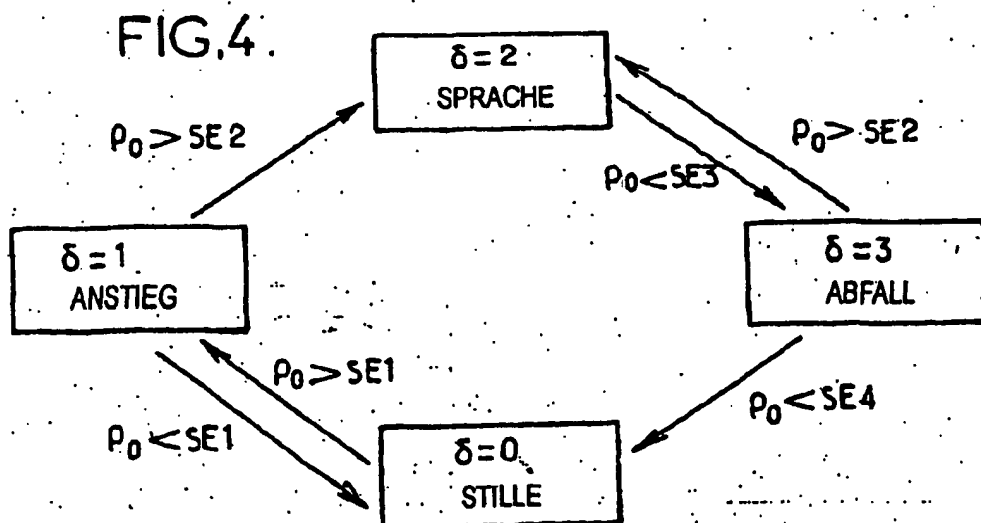
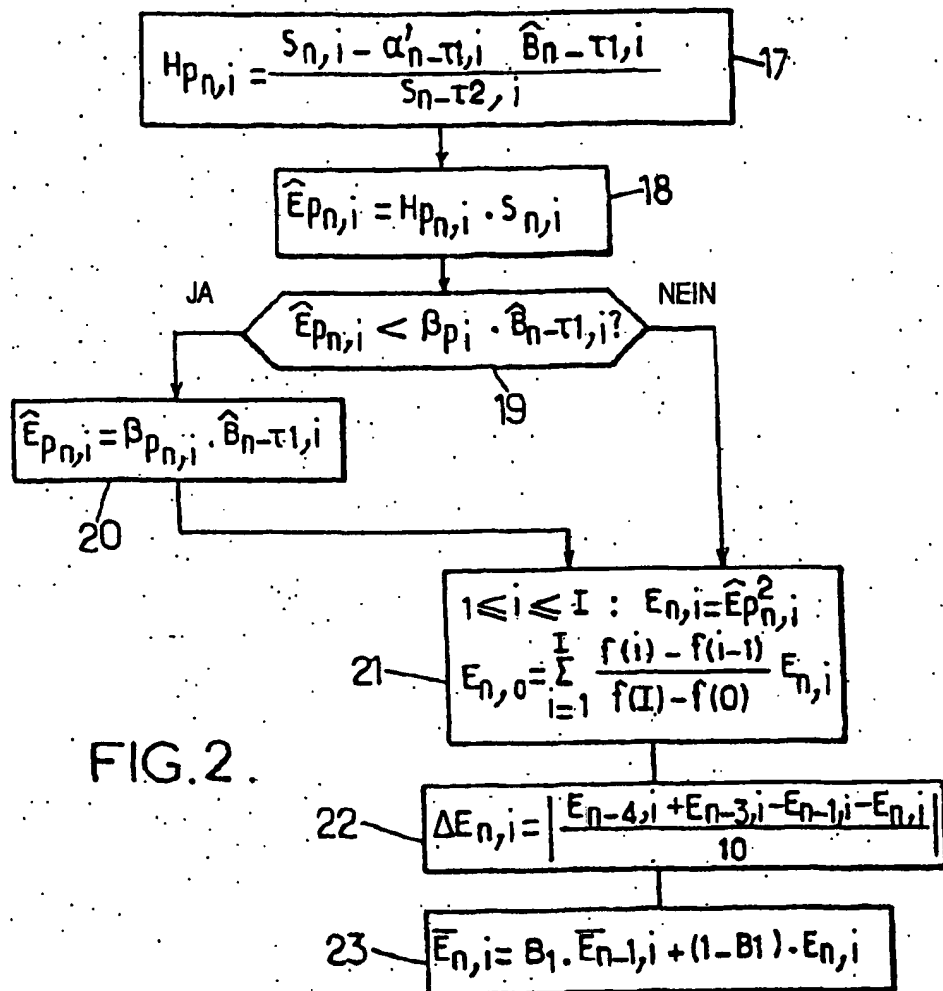
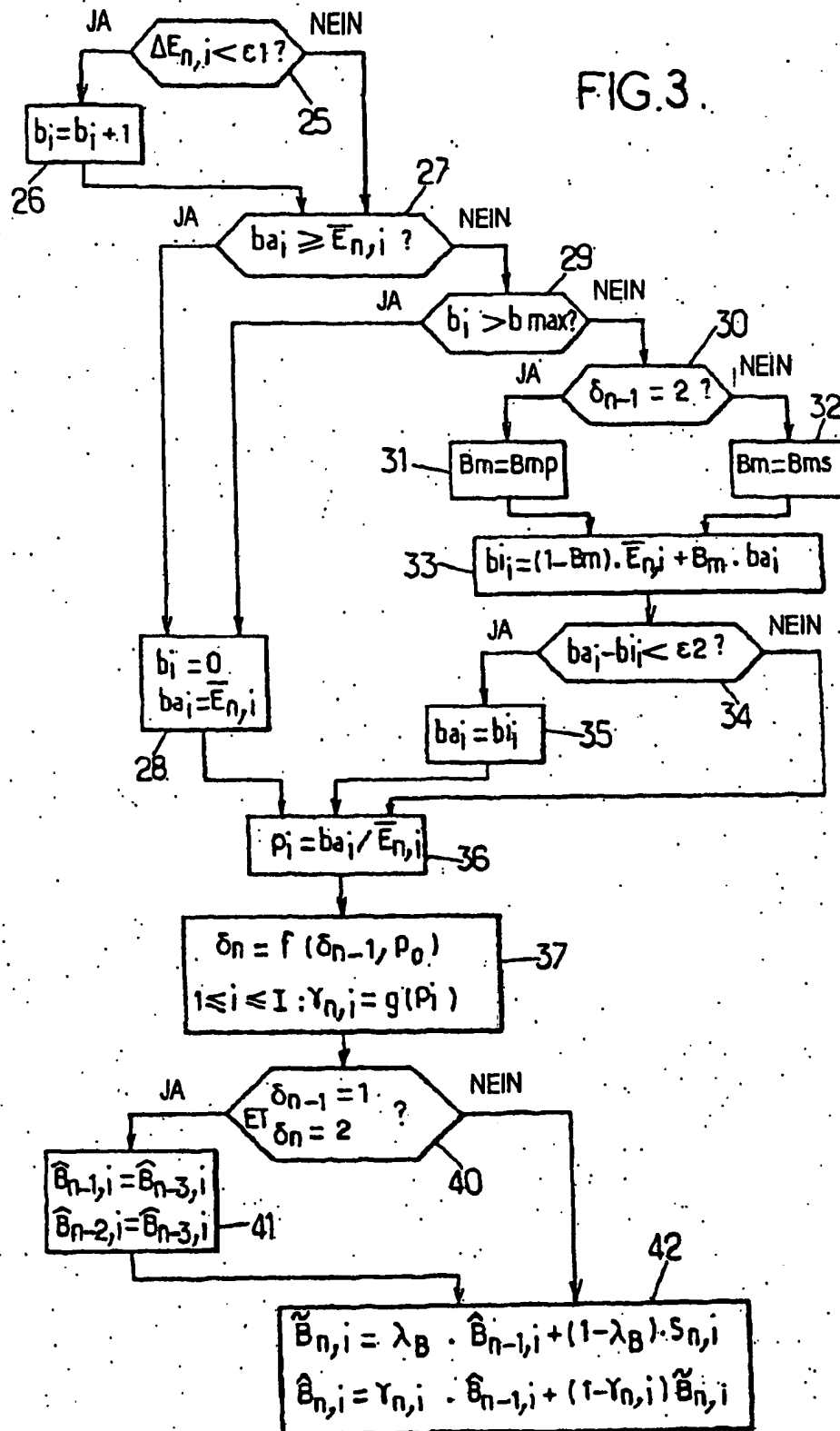


FIG.3.



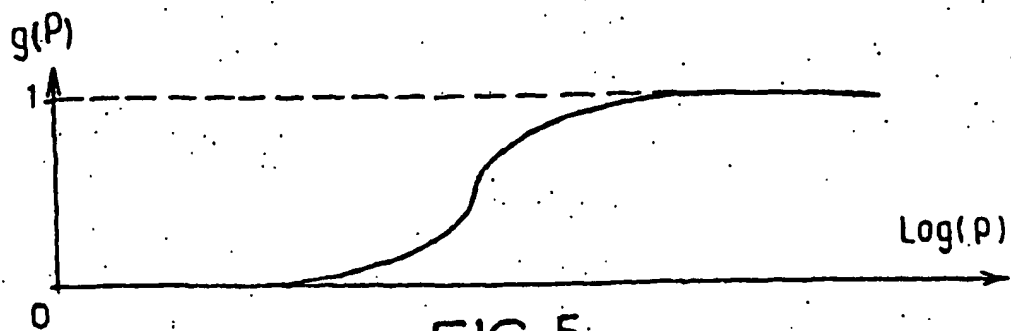


FIG. 5.

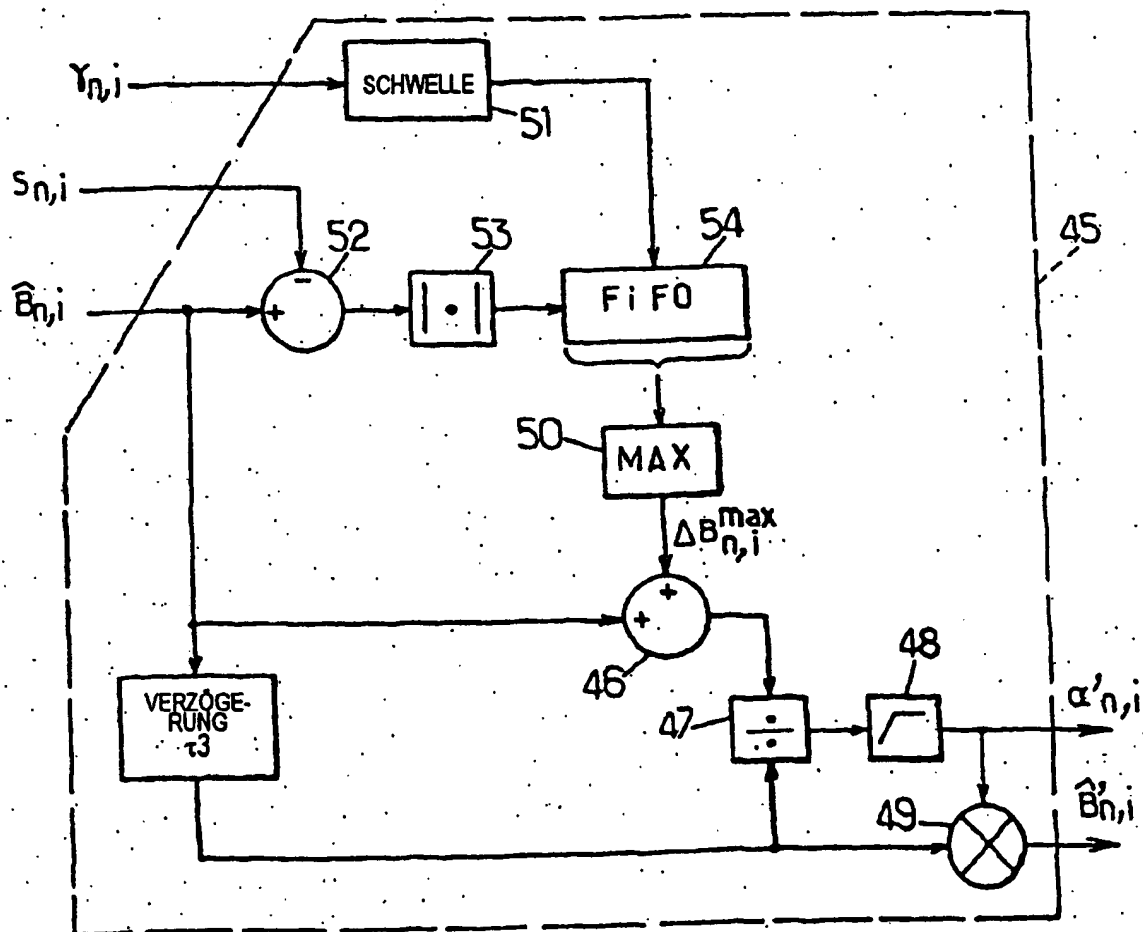


FIG. 6.

FIG. 7.

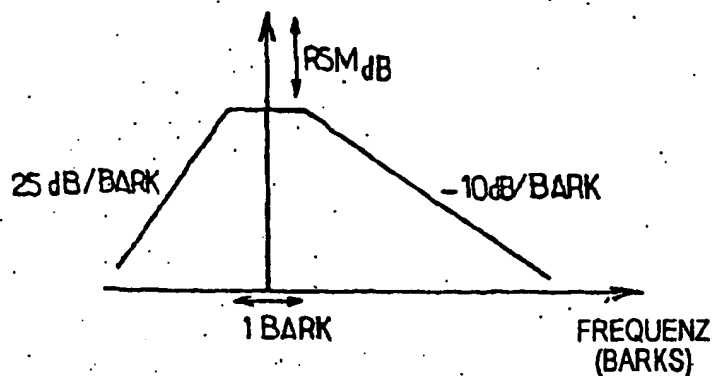
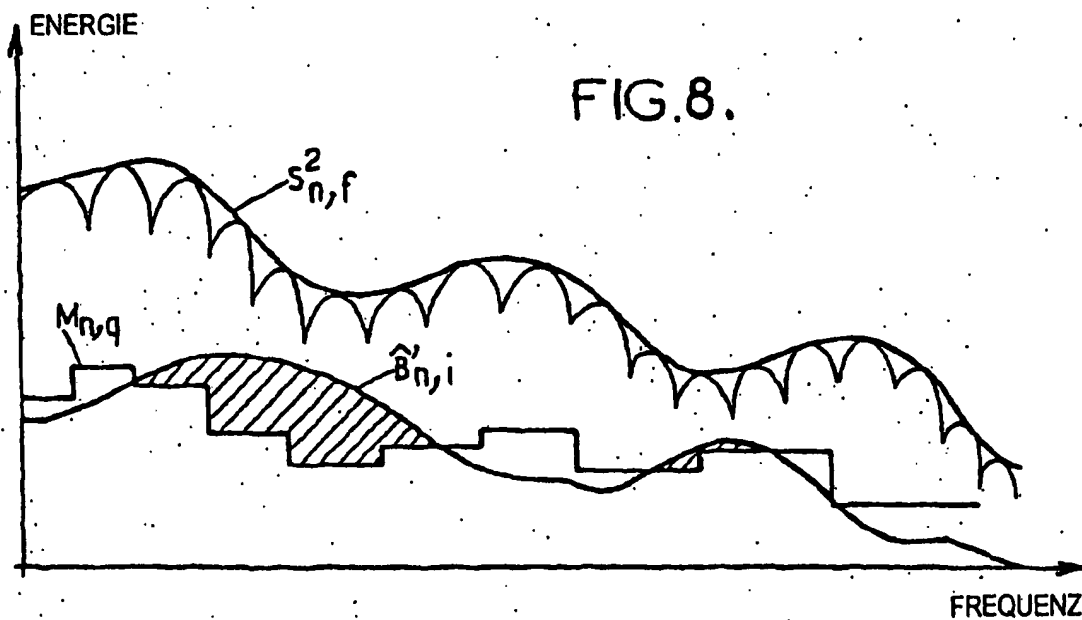


FIG. 8.



**This Page is Inserted by IFW Indexing and Scanning
Operations and is not part of the Official Record**

BEST AVAILABLE IMAGES

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images include but are not limited to the items checked:

- ☐ BLACK BORDERS
- ☐ IMAGE CUT OFF AT TOP, BOTTOM OR SIDES
- ☐ FADED TEXT OR DRAWING
- ☒ BLURRED OR ILLEGIBLE TEXT OR DRAWING
- ☐ SKEWED/SLANTED IMAGES
- ☐ COLOR OR BLACK AND WHITE PHOTOGRAPHS
- ☐ GRAY SCALE DOCUMENTS
- ☐ LINES OR MARKS ON ORIGINAL DOCUMENT
- ☐ REFERENCE(S) OR EXHIBIT(S) SUBMITTED ARE POOR QUALITY
- ☐ OTHER: _____

IMAGES ARE BEST AVAILABLE COPY.

As rescanning these documents will not correct the image problems checked, please do not report these problems to the IFW Image Problem Mailbox.